



DD2434 – Advanced Machine Learning

Lecture 2: **Probabilistic approach – fundamentals**

Pawel Herman

Computational Science and Technology (CST)

KTH Royal Institute of Technology

- Introduction
- Probabilistic approach
- Probability basics

Course content

- Aims and scope of this part of the course

I. Probabilistic regression

- *Lecture 2:* Probabilistic approach (objects and inference)
- *Lecture 3:* Linear regression
- *Lecture 4:* Kernels and introduction to Gaussian processes
- *Lecture 5:* Gaussian processes

II. Probabilistic representation learning

- *Lecture 6:* Representation learning

- **Introduction**
- Probabilistic approach
- Probability basics

Course content

- Learning activities
 - 5 lectures
 - 2 exercise sessions (fun with Gaussians ;-))
 - 1 assignment
 - project based on a research paper

- **Introduction**
- Probabilistic approach
- Probability basics

Short outline for today

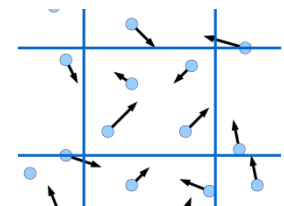
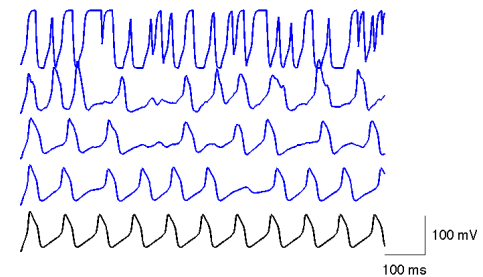
1. Introduction – why is probability theory relevant in ML?
2. Probabilistic approach to ML – principles (*recap*).
3. Probability basics (tbc).

- Introduction
- Probabilistic approach
- Probability basics

Introduction

Probabilistic approach

- Ubiquitous nature of uncertainty
 - imprecision, noise in data,
 - errors , missing information/data
 - gaps in knowledge, simplified description



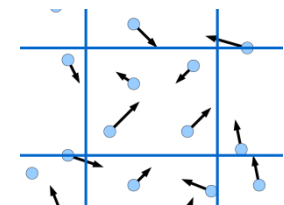
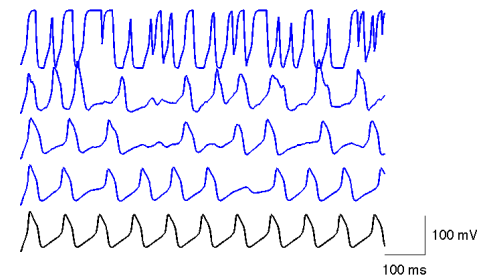
- Introduction
- Probabilistic approach
- Probability basics

Introduction

Probabilistic approach

- Ubiquitous nature of uncertainty
 - imprecision, noise in data,
 - errors , missing information/data
 - gaps in knowledge, simplified description

“The theory of probabilities is at bottom nothing but common sense reduced to calculus; it enables us to appreciate with exactness that which accurate minds feel with a sort of instinct for which of times they are unable to account.” (Laplace)

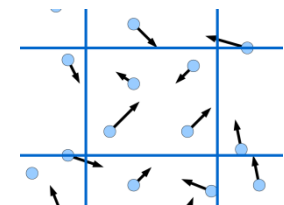
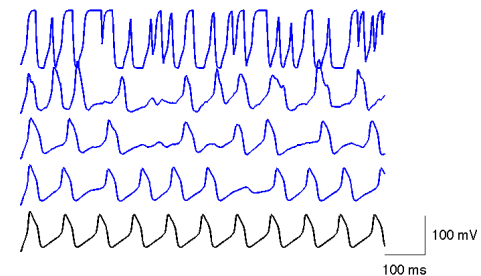


- Introduction
- Probabilistic approach
- Probability basics

Introduction

Probabilistic approach

- Ubiquitous nature of uncertainty
 - imprecision, noise in data
 - errors , missing information/data
 - gaps in knowledge, simplified description
- Probability theory provides a framework for modelling, reasoning etc. under uncertainty
 - unified, universal, intuitive, interpretable
 - beyond randomness, it is about uncertainty!
 - p. distributions as carriers or information (Jaynes, 2003)



Probabilistic perspective in ML

- Statistical ML: constructing stochastic models
 - fully probabilistic description and inference
 - theoretical assumptions, mathematical tractability, rigour
 - parameters estimated from observed data (learning)
 - interpretability and extra insights
 - machinery to propagate and account for uncertainty effects

Probabilistic perspective in ML

- Statistical ML: constructing stochastic models
 - fully probabilistic description and inference
 - theoretical assumptions, mathematical tractability, rigour
 - parameters estimated from observed data (learning)
 - interpretability and extra insights
 - machinery to propagate and account for uncertainty effects

BUT: *can be very hard for large-scale problems* and
difficult to derive solutions in a closed form

Probabilistic perspective in ML

- Statistical ML: constructing stochastic models

- fully probabilistic description and inference

- t

- p

- il

- n

“(statistical ML) provides the basis for learning algorithms that directly manipulate probabilities, as well as a framework for analyzing the operation of other algorithms that do not explicitly manipulate probabilities.”

Mitchell

rigour

effects

Probabilistic perspective in ML

- Philosophy of Bayesian approach
 - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)

Probabilistic perspective in ML

- Philosophy of Bayesian approach
 - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
 - Apply Bayesian machinery to propagate uncertainty
 - product and sum probability rules, Bayesian theorem

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

Probabilistic perspective in ML

- Philosophy of Bayesian approach
 - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
 - Apply Bayesian machinery to propagate uncertainty
 - product and sum probability rules, Bayesian theorem
 - the power of marginalisation

$$p(x_1, x_2, \dots, x_{n-1}) = \int p(x_1, x_2, \dots, x_n) dx_n$$

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

Probabilistic perspective in ML

- Philosophy of Bayesian approach
 - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
 - Apply Bayesian machinery to propagate uncertainty
 - Combine uncertain knowledge with data to reduce uncertainty (based on evidence from observations)
 - batch or sequence where posterior is iteratively updated

Probabilistic perspective in ML

- Philosophy of Bayesian approach
 - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
 - Apply Bayesian machinery to propagate uncertainty
 - Combine uncertain knowledge with data to reduce uncertainty (based on evidence from observations)
 - Two levels of inference:
parameter estimation and model selection (see Lecture 3)

Other ML approaches

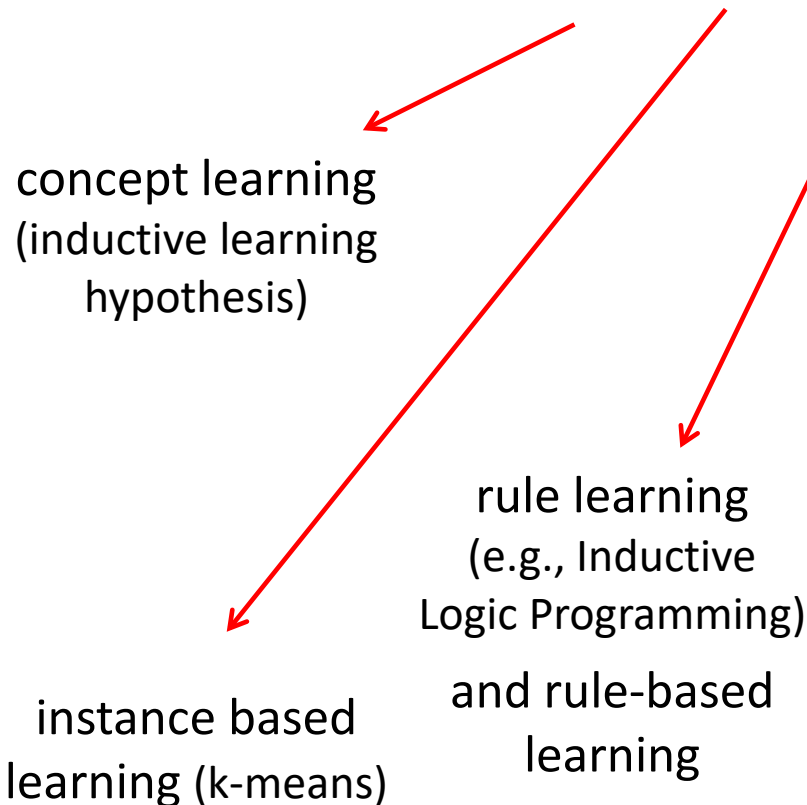
- Probabilistic narration is not the only one in ML



concept learning
based on symbolic
representations
(inductive learning
hypothesis)

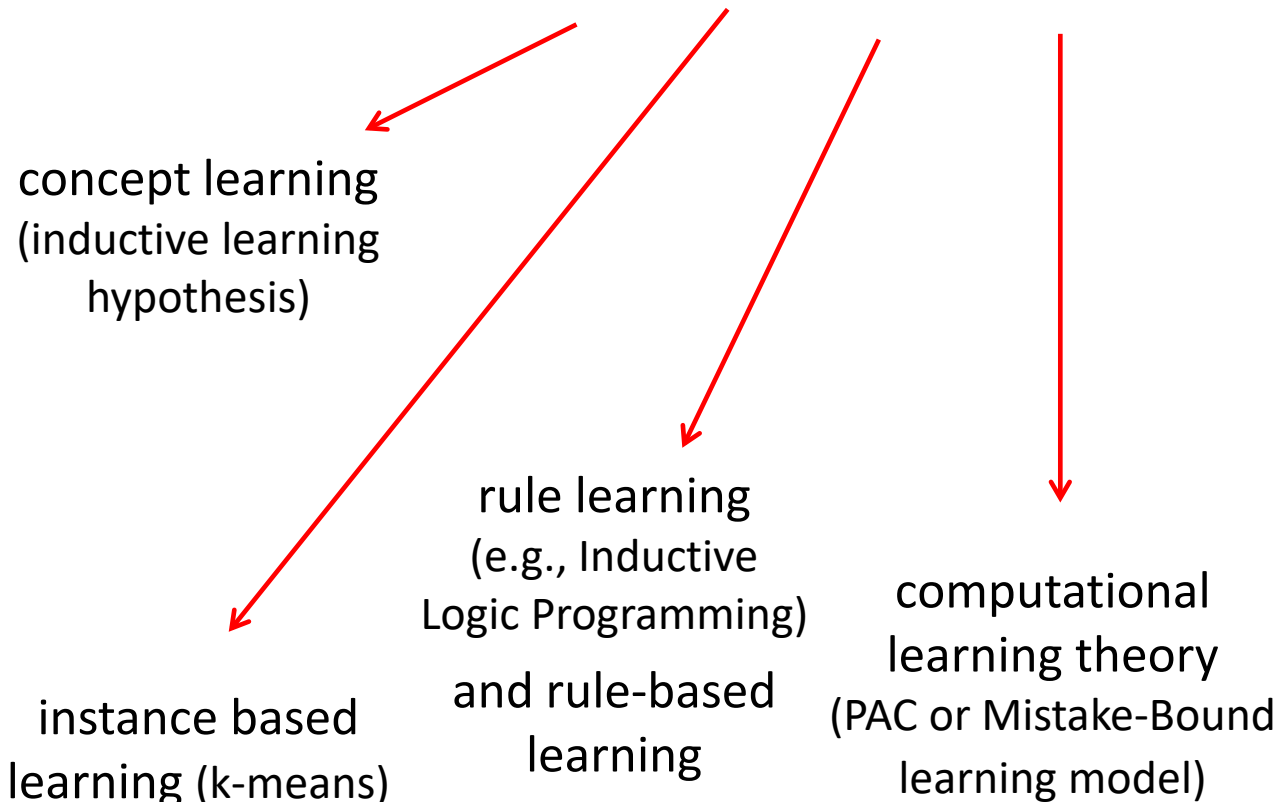
Other ML approaches

- Probabilistic narration is not the only one in ML



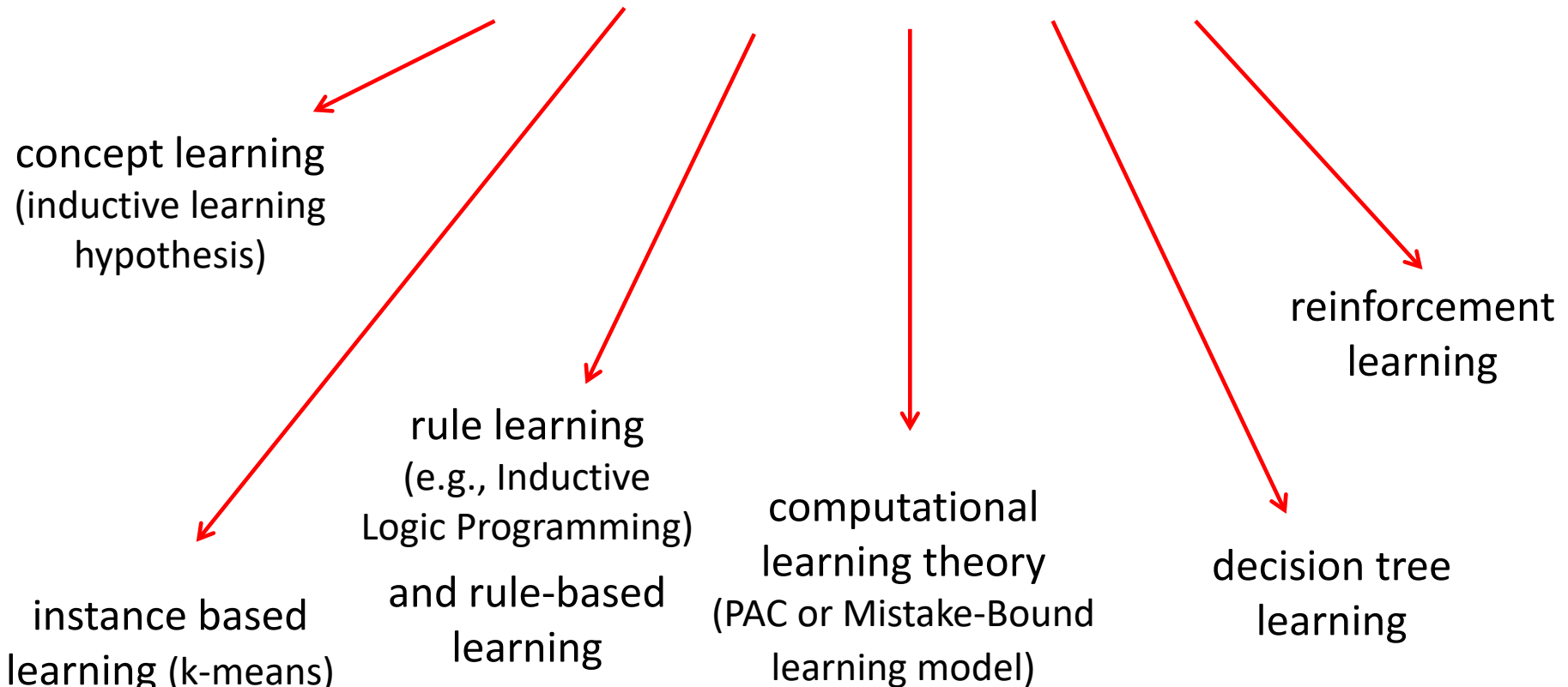
Other ML approaches

- Probabilistic narration is not the only one in ML



Other ML approaches

- Probabilistic narration is not the only one in ML



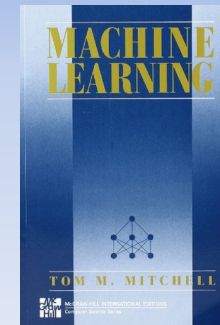
Other ML approaches

- Probabilistic narration is not the only one in ML

concept learning
based on symbolic
representations
(inductive learning)
hypothesis

You also may have heard about

- artificial neural networks, deep architectures
- search heuristics, e.g. genetic algorithms
- support vector machines
(origins in comp learning theory)
- etc.



reinforcement
learning

instance based
learning (k-means)

and rule-based
learning

(PAC or Mistake-Bound
learning model)

decision tree
learning

Other ML approaches

- Probabilistic narration is not the only one in ML

You also may have heard about

- artificial neural networks, deep architectures
- search heuristics, e.g. genetic algorithms
- support vector machines (origins in comp learning theory)
- etc.

ML comes in many flavours

concept learning
based on symbolic
representations
(inductive learning)
hypothesis

instance based
learning (k-means)

and rule-based
learning

(PAC or Mistake-Bound
learning model)

decision tree
learning

reinforcement
learning

Back to the probabilistic mindset

Bayesian vs frequentist perspective

I. Probability is a measure of belief (plausibility)

I. The ratio of outcomes in repeated trials.

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|--|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|---|--|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from data (likely outcome of exp.) |

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|---|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated. | IV. Each repeated experiment starts from ignorance. |

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|---|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated. | IV. Each repeated experiment starts from ignorance. |
| V. Estimators are good for available data. | V. Estimators are averaged across many trials. |

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|--|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated | IV. Each repeated experiment starts from ignorance. |
| V. Estimators are good for available data. | V. Estimators are averaged across many trials. |
| VI. Probability of a hypothesis given the data (posterior distribution). | VI. Probability of the data given hypothesis (likelihood, sampling dist.). |

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|---|--|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random realisation. |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated | IV. Each repeated experiment starts from ignorance. |
| V. Estimators are good for available data. | V. Estimators are averaged across many trials. |
| VI. Probability of a hypothesis given the data. | VI. Probability of the data given hypothesis. |
| VII. All variables/parameters have distribution. | VII. Parameters are fixed unknowns that can be point estimated from repeated trials. |

- Introduction
- **Probabilistic approach**
- Probability basics

- **ML – probabilistic perspective**
- Learning and inference
- Model complexity

Back to the probabilistic mindset

Bayesian vs frequentist perspective

- | | |
|--|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
| II. The data is fixed, models have probabilities | II. There is a true model and the data is a random |
| III. There does not
declaring proba | III. The data is produced from
one of exp.) |
| IV. Can incorporate prior knowledge,
probabilities can be updated | IV. Each repeated experiment starts from
ignorance. |
| V. Estimators are good for available data. | V. Estimators are averaged across many trials. |
| VI. Probability of a hypothesis given the data. | VI. Probability of the data given hypothesis. |
| VII. All variables/parameters have distribution. | VII. Parameters are fixed unknowns that can be
point estimated from repeated trials. |

Why isn't everyone Bayesian? (*Efron, 1986*)

Learning as inference – recap (c.f. DD2431, L6 G. Salvi)

- Learning distributions

➤ *“integrating data with domain knowledge of model environment”*

$$y = f(x) : X \rightarrow Y$$

$$\{(x_1, t_1), \dots, (x_N, t_N)\}$$

Learning implies the parameter identification of the model adopted to account for the mapping $X \rightarrow Y$

Inference is about using the model to predict the output corresponding to the sample input: *what is the model's response for $x = x^*$?*

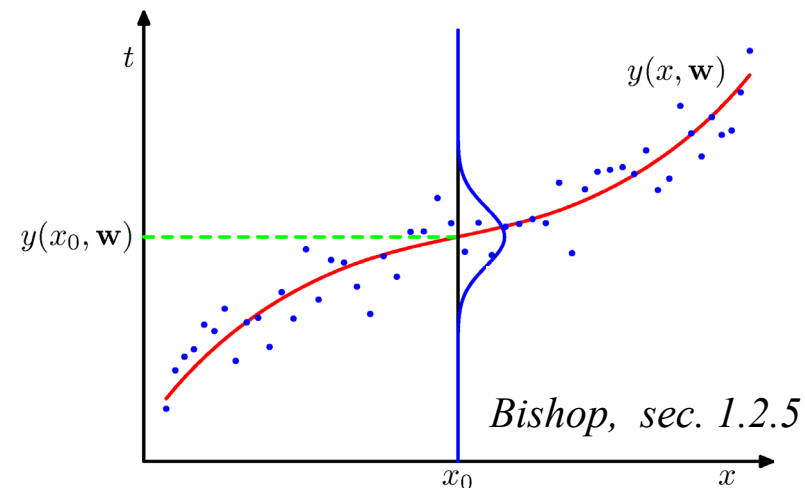
Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$



Learning as inference – recap

- Learning distributions

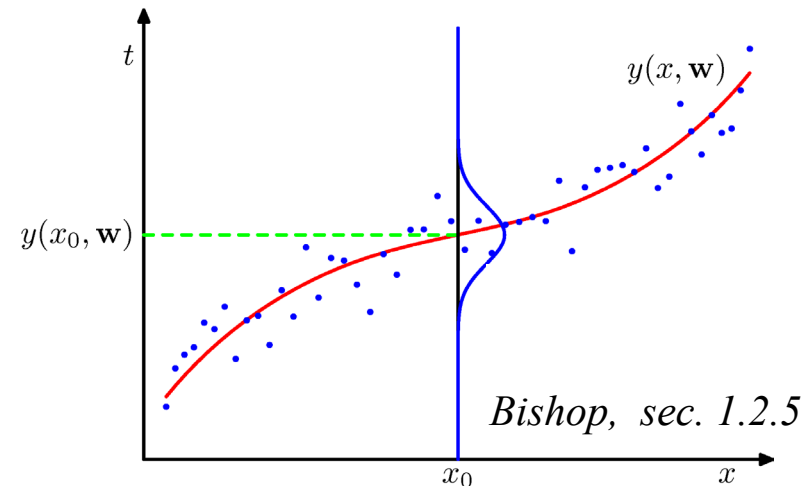
➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

Remarks about notation:

- 1) here we deal with one-dim input, x and output, t
- 2) parameters \mathbf{w} still constitute a vector (e.g. could be polynomial coefficients)
- 3) \mathbf{x} and \mathbf{t} , refer to the collection of all inputs and outputs, conceptually corresponding to D_x and D_t but in the vector form, so $\mathbf{x} \rightarrow \mathbf{t}$.



Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

➤ probabilistic framework

$p(\mathbf{w})$ – a prior probability distribution

$p(\mathcal{D} | \mathbf{w})$ – the likelihood function for $\mathcal{D} = \mathbf{t}$
(not a probability distrib. over \mathbf{w} , just a conditional probability)

$p(\mathbf{w} | \mathcal{D})$ – a posterior probability

Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

observation $\mathbf{t} = y(\mathbf{x}, \mathbf{w}) + \epsilon$

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

probab $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)^T$

$$p(\mathbf{w}) \quad \text{posterior} \propto \text{likelihood} \times \text{prior}$$

$p(\mathcal{D} | \mathbf{w})$ – the likelihood function for $\mathcal{D} = \mathbf{t}$
(not a probability distrib. over \mathbf{w} , just a conditional probability)

$p(\mathbf{w} | \mathcal{D})$ – a posterior probability

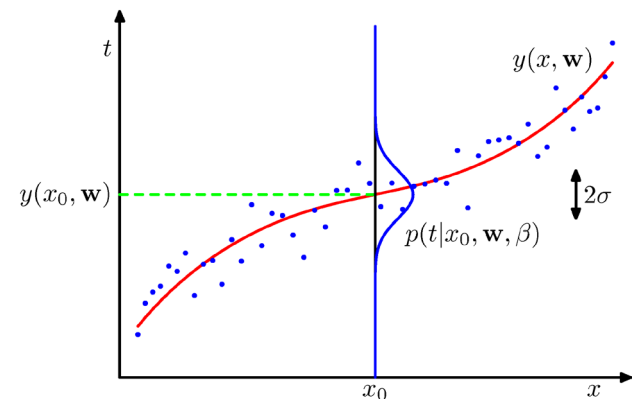
Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$



Bishop, sec. 1.2.5

Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

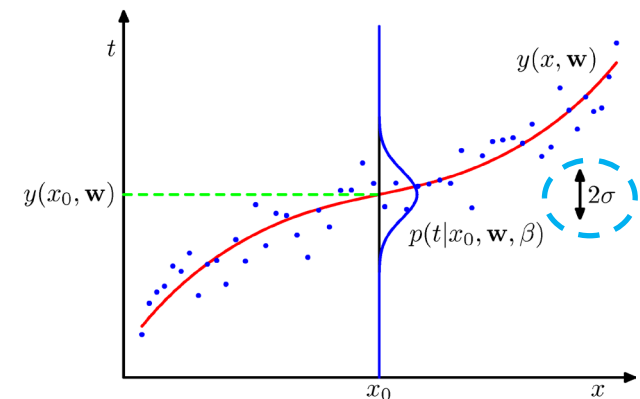
observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

Uncertainty (noise) in target data:

$$p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$$



precision



Bishop, sec. 1.2.5

Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

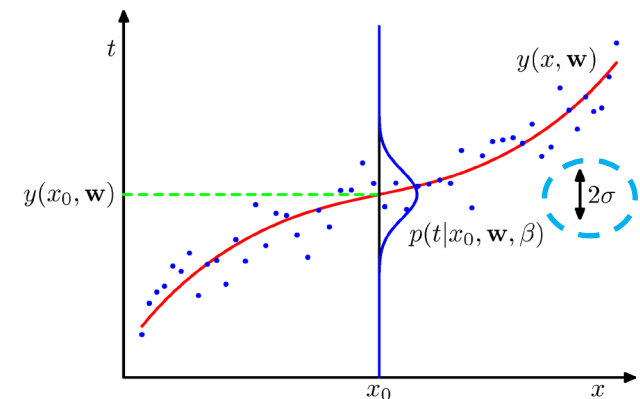
observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

Uncertainty (noise) in target data:

$$\underbrace{p(t | x, \mathbf{w}, \beta)}_{\text{predictive distribution}} = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$$

predictive distribution
(not a point estimate)

↑
precision



Bishop, sec. 1.2.5

Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}): X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations $\{(x_i, t_i): i=1, \dots, N\}$: $\mathbf{x} = (x_1, \dots, x_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$

SOLUTION 1

We want to find parameters to be able to use predictive distribution, $p(t|x)$, and infer the target:

$$\mathbb{E}[t | x] = \int t p(t | x, \mathbf{w}_{\text{opt}}, \beta_{\text{opt}}) dt$$

Learning as inference – recap

So, we follow the **max likelihood** approach

ML function for t (i.i.d.) under Gaussian noise $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

Learning as inference – recap

So, we follow the **max likelihood** approach

ML function for t (i.i.d.) under Gaussian noise $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

The log-likelihood:

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Learning as inference – recap

So, we follow the **max likelihood** approach

ML function for t (i.i.d.) under Gaussian noise $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

The log-likelihood:

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Still, input x is one-dim and \mathbf{x} & \mathbf{t} refer to the collection of all data points.

Learning as inference – recap

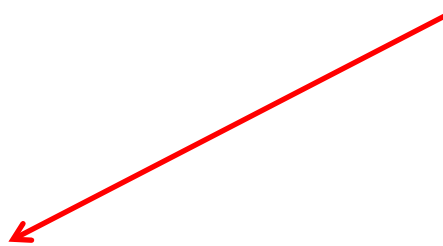
Maximum likelihood (ML) estimate:

Maximise $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$

Learning as inference – recap

Maximum likelihood (ML) estimate:

Maximise $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$


$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$

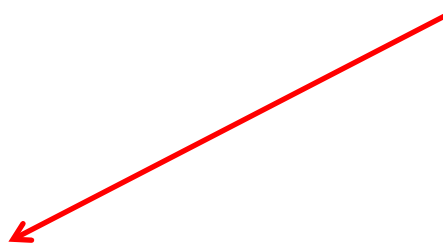
the sum-of-squares error function

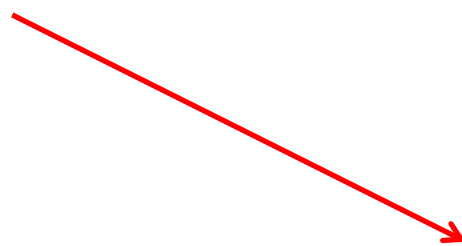
under the assumption of Gaussian noise : $p(t \mid x, \mathbf{w}, \beta) = \mathcal{N}(t \mid y(x, \mathbf{w}), \beta^{-1})$

Learning as inference – recap

Maximum likelihood (ML) estimate:

Maximise $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$


$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$


$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$

Learning as inference – recap

Maximum likelihood (ML) estimate:

Maximise $\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$

$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$



$$t_{\text{out}} = \mathbb{E}[t | x_{\text{in}}] = \int t p(t | x_{\text{in}}, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) dt$$



Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}) : X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

SOLUTION 2

Parameters of the model could also be random variables with a prior distribution, $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$

Now, we must **maximise the posterior** $p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta)$, i.e. find the most probable \mathbf{w} given data:

$$\max p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

Learning as inference – recap

- Learning distributions

➤ curve-fitting example: $y(x, \mathbf{w}) : X \rightarrow Y$ (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

SOLUTION 2

Parameters of the model could also be random variables with a prior distribution, $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$

Now, we must **maximise the posterior** $p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta)$, i.e. find the most probable \mathbf{w} given data:

$$\max p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

Again, \mathbf{x} & \mathbf{t} refer to the data collection, not individual input or outputs

Learning as inference – recap

Maximum posterior (MAP) estimate:

Maximise $\ln p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta)$



$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$

Learning as inference – recap

Maximum posterior (MAP) estimate:

Maximise $\ln p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta)$



$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$

BUT: SOLUTIONS 1 (ML) and 2 (MAP) give point estimates of \mathbf{w}

Learning as inference – Bayesian view

Instead of estimating “optimal” parameters \mathbf{w} , let’s integrate over all values of \mathbf{w} (let’s make use of the distribution)

Marginalisation:

$$p(t \mid x, \mathbf{x}, \mathbf{t}, \alpha, \beta) = \int p(t \mid x, \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta) d\mathbf{w}$$

predictive “noise” model posterior

We assume we know what α and β are.

Learning as inference – Bayesian view

Instead of estimating “optimal” parameters \mathbf{w} , let's integrate over all values of \mathbf{w} for our linear model $y = \sum_{i=0}^M w_i \phi_i(x)$

For a one-dim polynomial model: $\phi_i(x) = x^i$ (order $M-1$)

Learning as inference – Bayesian view

Instead of estimating “optimal” parameters \mathbf{w} , let’s integrate over all values of \mathbf{w} for our linear model $y = \sum_{i=0}^M w_i \phi_i(x)$

For a one-dim polynomial model: $\phi_i(x) = x^i$

$$p(t \mid x, \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t \mid m(x), s^2(x)) \begin{cases} m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S} \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \\ s^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S} \boldsymbol{\phi}(x) \end{cases}$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

$$\boldsymbol{\phi}(x) = (\phi_0(x), \dots, \phi_M(x))$$

Learning as inference – summary

In summary:

To use a predictive distribution and infer the output t for the given input x_{in}

Learning as inference – summary

In summary:

To use a predictive distribution and infer the output t for the given input x_{in}

.....ML and MAP approaches produce point estimates of \mathbf{w}

$$\text{ML:} \quad \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP:} \quad \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

Learning as inference – summary

In summary:

To use a predictive distribution and infer the output t for the given input x_{in}

.....ML and MAP approaches produce point estimates of \mathbf{w}

$$\text{ML:} \quad \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP:} \quad \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

.....in Bayesian approach \mathbf{w} is integrated over (**marginalisation**)

$$\text{Bayes:} \quad \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$

c.f. DD2431, Lecture 6, G. Salvi

Learning as inference – summary

In summary:

To use a predictive distribution and infer the output t for the given input x_{in}

.....ML and MAP approaches produce point estimates of \mathbf{w}

$$\text{ML: } \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP: } \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$



Use $p(t | x_{in}, \mathbf{w}_{\text{ML}})$ or $p(t | x_{in}, \mathbf{w}_{\text{MAP}})$ for prediction.

.....in Bayesian approach \mathbf{w} is integrated over (**marginalisation**)

$$\text{Bayes: } \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$



Marginalise over \mathbf{w} :

$$p(t | \mathcal{D}) = \int p(t | x_{in}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

Learning as inference – summary

In summary:

To use a predictive distribution and infer the output t for the given input x_{in}

.....ML and MAP approaches produce point estimates of \mathbf{w}

$$\text{ML: } \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP: } \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

Frequentist philosophy

$p(t | x_{in}, \mathbf{w}_{\text{MAP}})$ for prediction.

.....in Bayesian approach \mathbf{w} is integrated over (**marginalisation**)

$$\text{Bayes: } \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$

Bayesian philosophy

$$p(t | \mathcal{D}) = \int p(t | x_{in}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

Inference and decision (classification)

The ***inference*** stage of classification $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

Model the inputs \mathbf{x} and outputs C

$$\left. \begin{array}{l} p(\mathbf{x} | C_k) \\ p(C_k) \end{array} \right\} \text{ for each class } C_k \quad p(\mathbf{x}, C_k)$$


$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k)$$

Inference and decision (classification)

The ***inference*** stage of classification $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

Model the inputs \mathbf{x} and outputs C

$$\left. \begin{array}{l} p(\mathbf{x} | C_k) \\ p(C_k) \end{array} \right\} \text{ for each class } C_k$$


$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k)$$

GENERATIVE approach

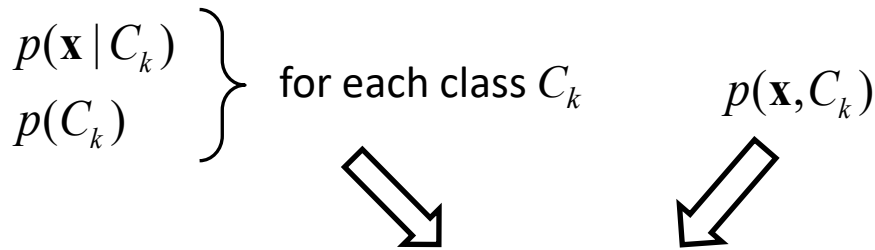
Remarks:

- 1) K classes
- 2) \mathbf{x} – multi-dim input feature vector

Inference and decision (classification)

The ***inference*** stage of classification $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

Model the inputs \mathbf{x} and outputs C



$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k)$$

GENERATIVE approach

Solve first the inference problem of determining **posteriors** for each class without modelling $p(C_k, \mathbf{x})$ or $p(\mathbf{x} | C_k)$

$$p(C_k | \mathbf{x})$$

DISCRIMINATIVE approach

Generative vs discriminative approach

What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)

Generative vs discriminative approach

What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)
- Easy and elegant way of handling missing or unlabelled data
- Generative model allows for..... generating data
 - > generative models can be run “backwards”

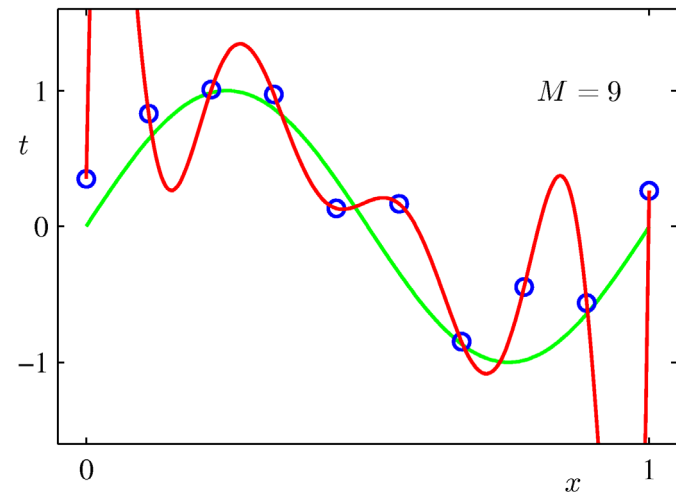
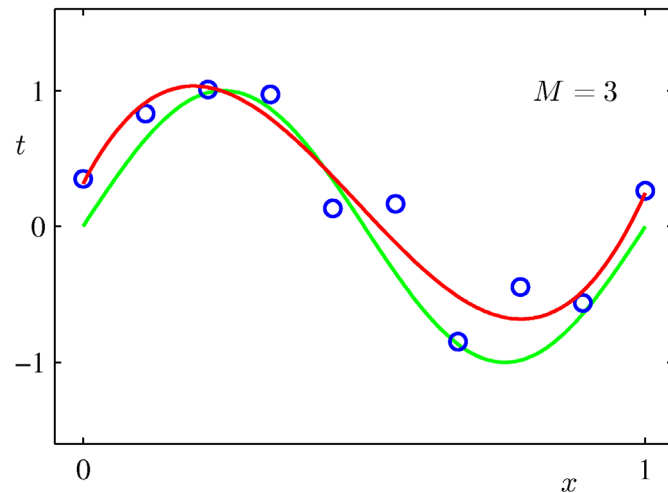
Generative vs discriminative approach

What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)
- Easy and elegant way of handling missing or unlabelled data
- Generative model allows for..... generating data
 - > generative models can be run “backwards”
- BUT: discriminative models tend to be more accurate (less vulnerable to assumptions)

Model complexity

- Overfitting

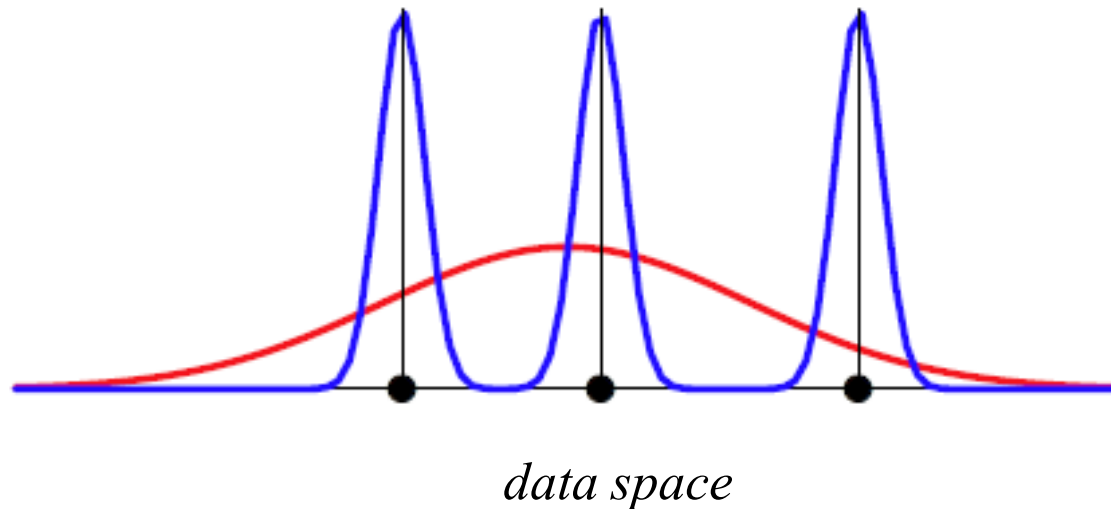


Bishop, sec. 1.1

Please recall from Lecture 1 (J. Lagergren)

Model complexity

- Overfitting



The **likelihood** can be arbitrarily large by adding parameters.

Please recall from Lecture 6, DD2431 (G. Salvi)

Frequentist viewpoint – bias-variance dilemma

- Maximum Likelihood estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$

Remarks about notation:

- 1) \mathbf{x} can be a multi-dim input
- 2) $y(\mathbf{x})$ really is $y(\mathbf{x}, \mathbf{w})$ that depends on parameters \mathbf{w}

Frequentist viewpoint – bias-variance dilemma

- **Maximum Likelihood** estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

Frequentist viewpoint – bias-variance dilemma

- **Maximum Likelihood** estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

$$\begin{aligned}\mathbb{E}[L] &= \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \\ &\quad \Downarrow \\ \mathbb{E}[L] &= \int \{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Frequentist viewpoint – bias-variance dilemma

- Maximum Likelihood estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \underbrace{\int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}}$$

noise on the data,
independent of the regressor $y(\mathbf{x})$

Frequentist viewpoint – bias-variance dilemma

- Maximum Likelihood estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

optimal prediction to be accounted for by
our regression model $y(\mathbf{x}): h(\mathbf{x}) = \mathbb{E}_t[t | \mathbf{x}]$

$$\mathbb{E}[L] = \int \underbrace{\{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2}_{\{y(\mathbf{x}) - h(\mathbf{x})\}^2} p(\mathbf{x}) d\mathbf{x} + \underbrace{\int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}}_{\text{noise on the data, independent of the regressor } y(\mathbf{x})}$$

Frequentist viewpoint – bias-variance dilemma

- Maximum Likelihood estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

optimal prediction to be accounted for by our regression model $y(\mathbf{x})$: $h(\mathbf{x}) = \mathbb{E}_t[t | \mathbf{x}]$

h is the underlying model (data generating mechanism) that we want to approximate.

$$\mathbb{E}[L] = \int \underbrace{\{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2}_{\{y(\mathbf{x}) - h(\mathbf{x})\}^2} p(\mathbf{x}) d\mathbf{x} + \underbrace{\int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}}_{\text{noise on the data, independent of the regressor } y(\mathbf{x})}$$

Frequentist viewpoint – bias-variance dilemma

- Maximum Likelihood estimate for linear regression corresponds to minimising loss $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$
- The expected loss, L , to be minimised looks then as follows:

$$\mathbb{E}[L] = \underbrace{\int \{y(\mathbf{x}) - \mathbb{E}_t[t | \mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x}}_{\text{has to be minimised}} + \underbrace{\int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}}_{\text{noise on the data, independent of the regressor } y(\mathbf{x})}$$

Frequentist approach

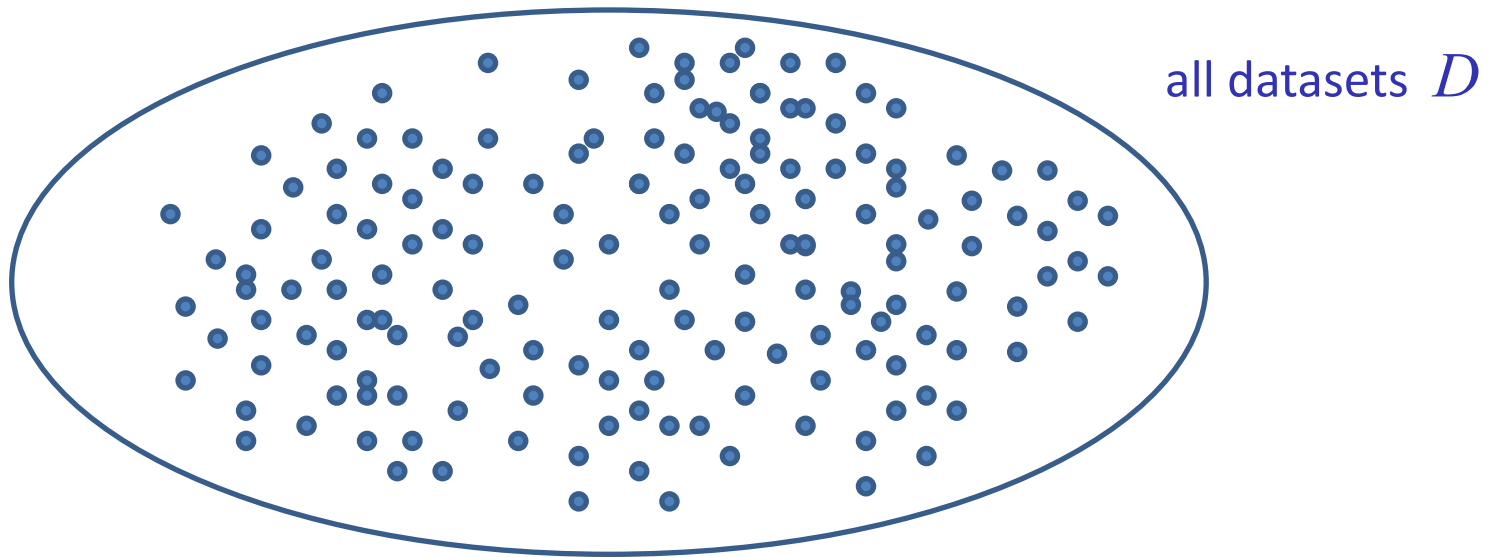
- Find optimal parameters for y using data D

$$\int_{\mathcal{D}} \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \rightarrow \min \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2$$

Frequentist approach

- Find optimal parameters for y using data D

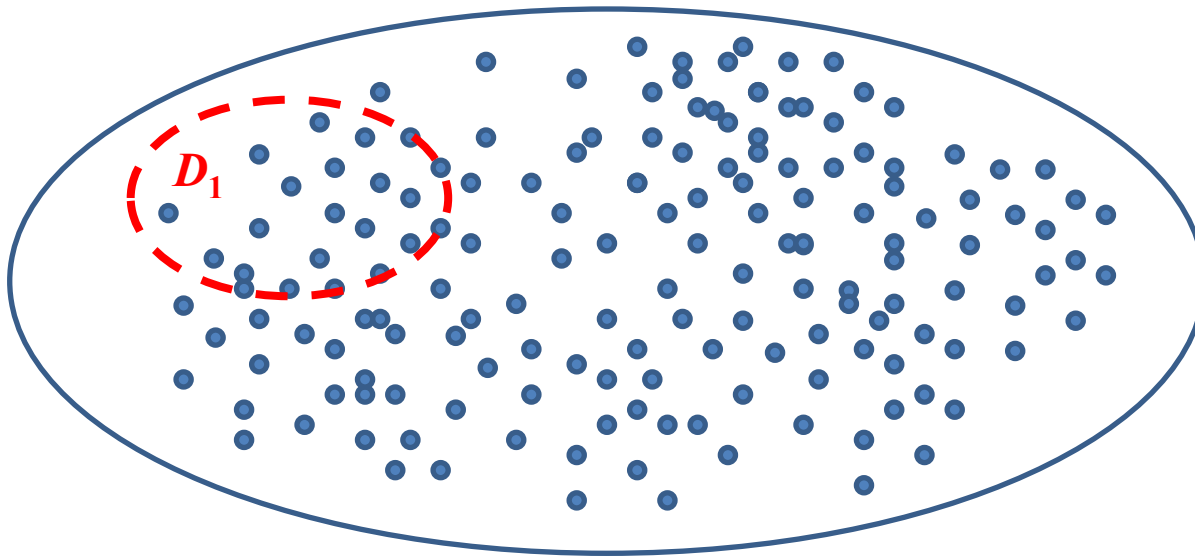
$$\int_{\mathcal{D}} \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \rightarrow \min \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2$$



Frequentist approach

- Find optimal parameters for y using data D

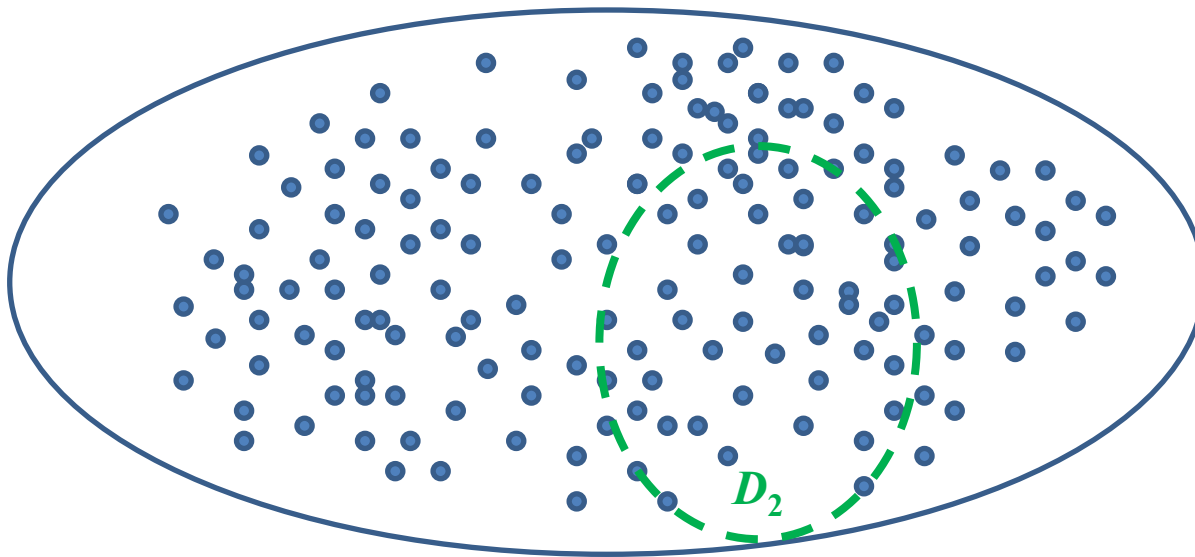
$$\min \{y(\mathbf{x}; \mathcal{D}_1) - h(\mathbf{x})\}^2$$



Frequentist approach

- Find optimal parameters for y using data D

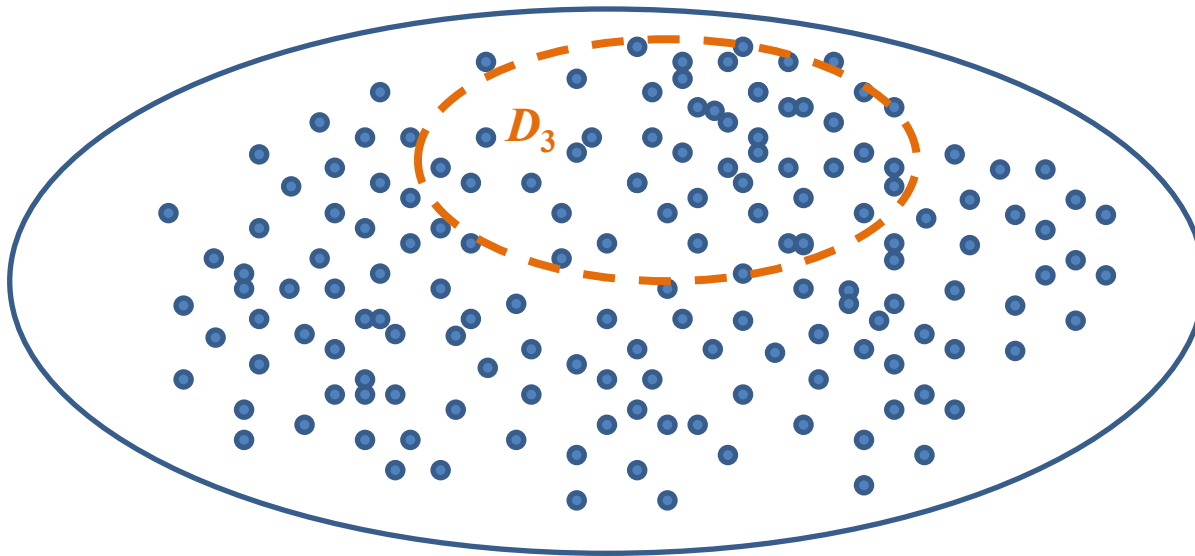
$$\min \{y(\mathbf{x}; \mathcal{D}_2) - h(\mathbf{x})\}^2$$



Frequentist approach

- Find optimal parameters for y using data D

$$\min \{y(\mathbf{x}; \mathcal{D}_3) - h(\mathbf{x})\}^2$$



Bias-variance decomposition

Expected value over all possible datasets D

$$\mathbb{E}[L] = \mathbb{E}_D \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

Bias-variance decomposition

Expected value over all possible datasets D

$$\mathbb{E}[L] = \mathbb{E}_D \left\{ y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

Taking frequentist perspective of fixing the model for specific data:

$$\left\{ y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right\}^2 = \left\{ y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_D [y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_D [y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right\}^2$$

Bias-variance decomposition

Expected value over all possible datasets D

$$\mathbb{E}[L] = \mathbb{E}_D \{y(\mathbf{x}; D) - h(\mathbf{x})\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

Taking frequentist perspective of fixing the model for specific data:

$$\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2 = \{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)] + \mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2$$

...and then averaging all possible data \mathbb{E}_D

$$\mathbb{E}_D \left[\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2 \right] = \{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 + \mathbb{E}_D \left[\{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}^2 \right]$$

Bias-variance decomposition

Expected value over all possible datasets D

$$\mathbb{E}[L] = \mathbb{E}_D \left\{ y(\mathbf{x}; D) - h(\mathbf{x}) \right\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

$$\left\{ y(\mathbf{x}; D) - h(\mathbf{x}) \right\}^2 = \left\{ y(\mathbf{x}; D) - \mathbb{E}_D [y(\mathbf{x}; D)] + \mathbb{E}_D [y(\mathbf{x}; D)] - h(\mathbf{x}) \right\}^2$$

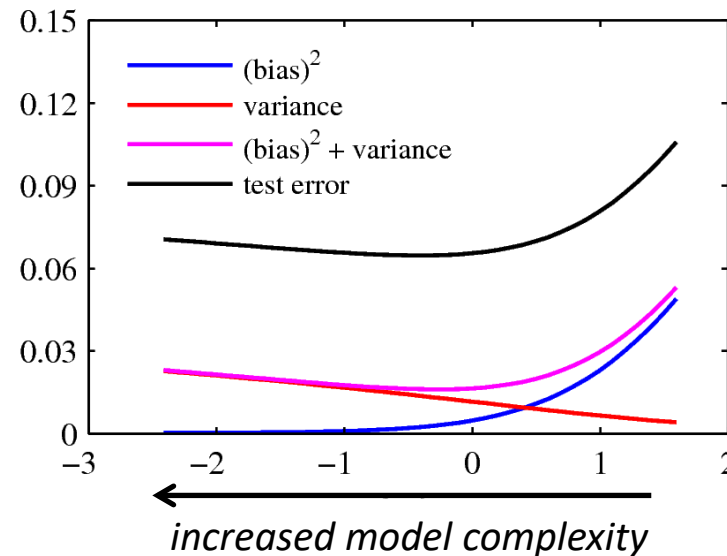
$$\mathbb{E}_D \left[\left\{ y(\mathbf{x}; D) - h(\mathbf{x}) \right\}^2 \right] = \underbrace{\left\{ \mathbb{E}_D [y(\mathbf{x}; D)] - h(\mathbf{x}) \right\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_D \left[\left\{ y(\mathbf{x}; D) - \mathbb{E}_D [y(\mathbf{x}; D)] \right\}^2 \right]}_{\text{variance}}$$

Bias-variance dilemma: model complexity

Expected value over all possible datasets D

$$\mathbb{E}_D[L] = \mathbb{E}_D \left\{ y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

$$\mathbb{E}[L] = (\text{bias})^2 + \text{variance} + \text{noise}$$



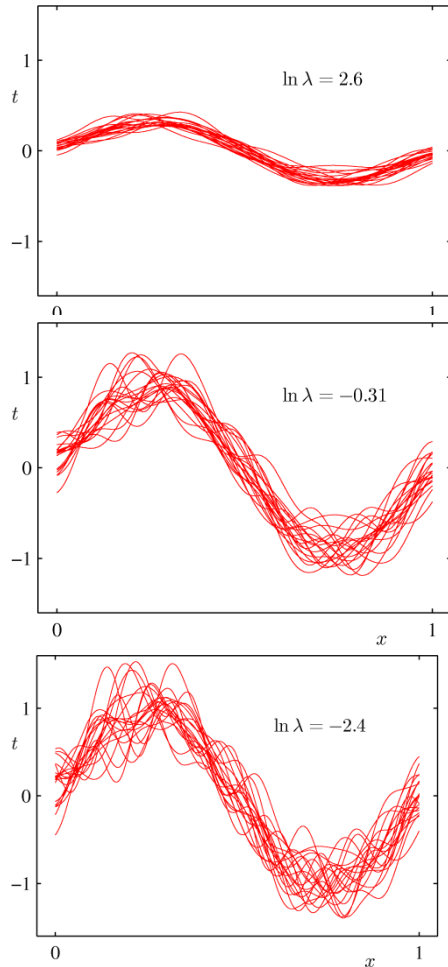
Bishop, sec. 3.2

- Introduction
- **Probabilistic approach**
- Probability basics

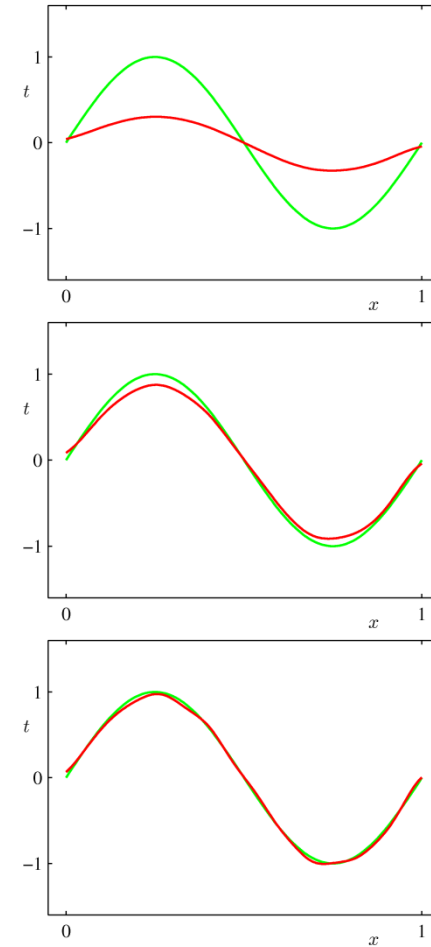
- ML – probabilistic perspective
- Learning and inference
- **Model complexity**

Bias-variance: illustrative example

Models for different data samples



increased model complexity



Model **average** and the **original**

Bishop, sec. 3.2

Model selection (*c.f.* L1, J. Lagergren)

- Occam's razor – *“Accept the simplest explanation that fits the data.”*

Model selection (*c.f.* L1, J. Lagergren)

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
 - bias-variance dilemma

Model selection (*c.f.* L1, J. Lagergren)

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
 - bias-variance dilemma
 - need to control the model's complexity
 - regularisation
 - correction for the bias of ML estimates (AIC, BIC)
 - empirical estimate of generalisation error on a hold-out set (validation, resampling)
 - structural risk minimization (SRM) (minimise upper bound on the true risk), see also VC dimension (statistical learning theory)

Model selection (*c.f.* L1, J. Lagergren)

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood

➤ bias

➤ need

Bayesian model selection (to appear in Lecture 3)

- • the second inference stage to compare models the Bayesian way
- • Bayes' factor and complexity penalisation (relies on the evidence approximation)
- structural risk minimization (SRM) (minimise upper bound on the true risk), see also VC dimension

- Introduction
- Probabilistic approach
- **Probability basics**

Selected fundamental concepts

1. Introduction – why is probability theory relevant in ML?
2. Probabilistic approach to ML – principles (*recap*).
- 3. Probability basics.**

Selected fundamental concepts

1. Introduction – why is probability theory relevant in ML?
2. Probabilistic approach to ML – principles (*recap*).
- 3. Probability basics.**

You should really be familiar with the fundamentals!

*Please take a closer look at **sec. 1.2, 2.1-2.4** in (Bishop, 2006)*

Here, we will VERY briefly mention

- the multivariate Gaussian distribution
- Central Limit Theorem

Selected fundamental concepts

1. Introduction – why is probability theory relevant in ML?
2. Probabilistic approach to ML – principles (*recap*).
- 3. Probability basics.**

You should really be familiar with the fundamentals!

Please

Here

More to come in
the first practical session

focus on sec.2.3- 2.4 (Bishop)

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian (normal) distribution

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{L/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$\mathbf{x} \in \mathbb{R}^L$

$\boldsymbol{\mu}$ – mean vector, \mathbb{R}^L

$\boldsymbol{\Sigma}$ – covariance matrix, $L \times L$

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian (normal) distribution

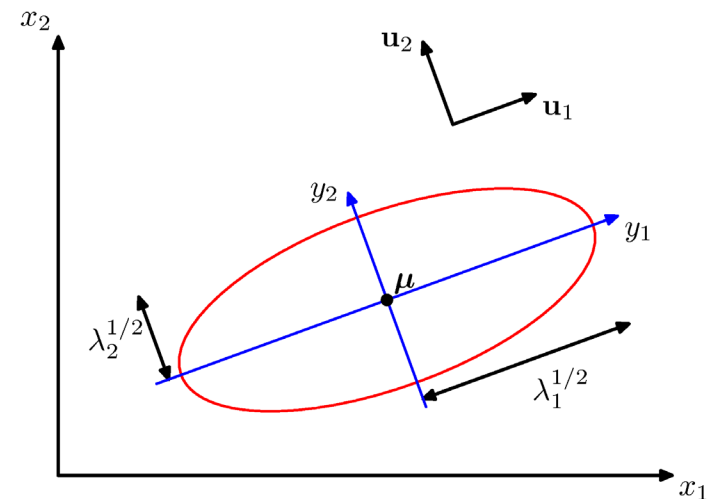
$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{L/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ \overbrace{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}^{\text{quadratic form}} \right\}$$

$\mathbf{x} \in \mathbb{R}^L$

$\boldsymbol{\mu}$ – mean vector, \mathbb{R}^L

$\boldsymbol{\Sigma}$ – covariance matrix, $L \times L$

$$\boldsymbol{\Sigma} = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$



Bishop, sec. 2.3

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian (normal) distribution

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{L/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ \overbrace{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}^{\text{quadratic form}} \right\}$$

$\mathbf{x} \in \mathbb{R}^L$

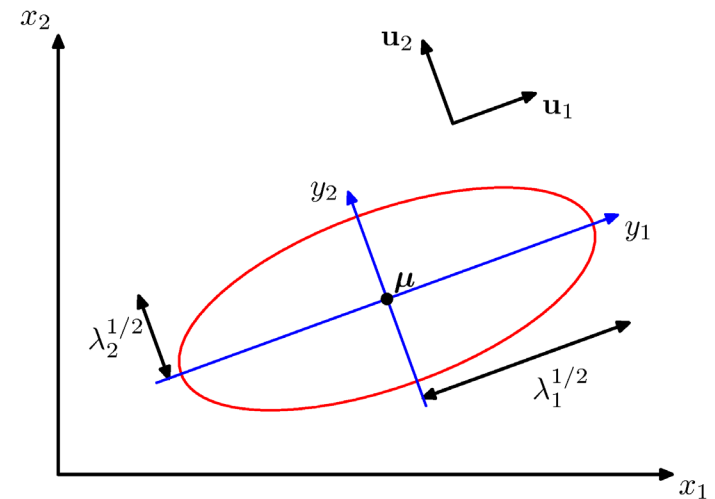
$\boldsymbol{\mu}$ – mean vector, \mathbb{R}^L

$\boldsymbol{\Sigma}$ – covariance matrix, $L \times L$

$$\mathbb{E}[\mathbf{x}] = \int \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbf{x} d\mathbf{x} = \boldsymbol{\mu}$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

$$\text{cov}[\mathbf{x}] = \mathbb{E}\left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T\right] = \boldsymbol{\Sigma}$$



Bishop, sec. 2.3

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian distribution

Two-dimensional (2D) case, i.e. $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

Remark about notation:

Unlike in the beginning, where x_i corresponded to the i -th realisation of one-dim r.v. x , **here** x_1 and x_2 correspond to one-dimensional variables that constitute two-dim r.v. \mathbf{x} .

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian distribution

Two-dimensional (2D) case, i.e. $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$$\mathbf{\Sigma} = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$\mu_i = \mathbb{E}[x_i]$$

$$\sigma_{i,j} = \sigma(x_i, x_j) = \mathbb{E} \left[(x_i - \mathbb{E}[x_i]) (x_j - \mathbb{E}[x_j])^T \right]$$

$$\sigma_{i,i} = \sigma_i^2 = \text{var}[x_i]$$

- Introduction
- Probabilistic approach
- **Probability basics**

The Gaussian distribution

Two-dimensional (2D) case, i.e. $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$\mu_i = \mathbb{E}[x_i]$$

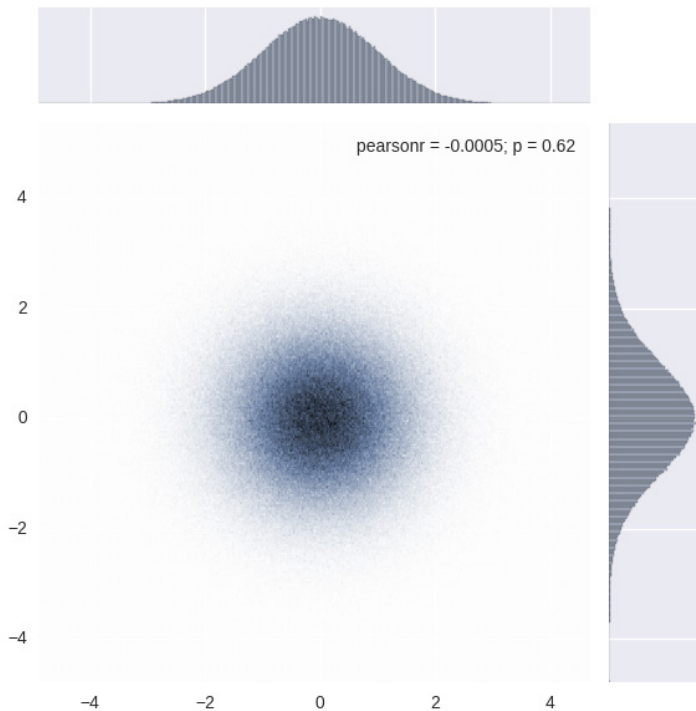
$$\sigma_{i,j} = \sigma(x_i, x_j) = \mathbb{E} \left[(x_i - \mathbb{E}[x_i]) (x_j - \mathbb{E}[x_j])^T \right]$$

$$\sigma_{i,i} = \sigma_i^2 = \text{var}[x_i]$$

Sample covariance:
$$\bar{\sigma}_{i,j} = \frac{1}{N-1} \sum_{k=1}^N (x_i^{(k)} - \mu_i) (x_j^{(k)} - \mu_j)$$

- Introduction
- Probabilistic approach
- **Probability basics**

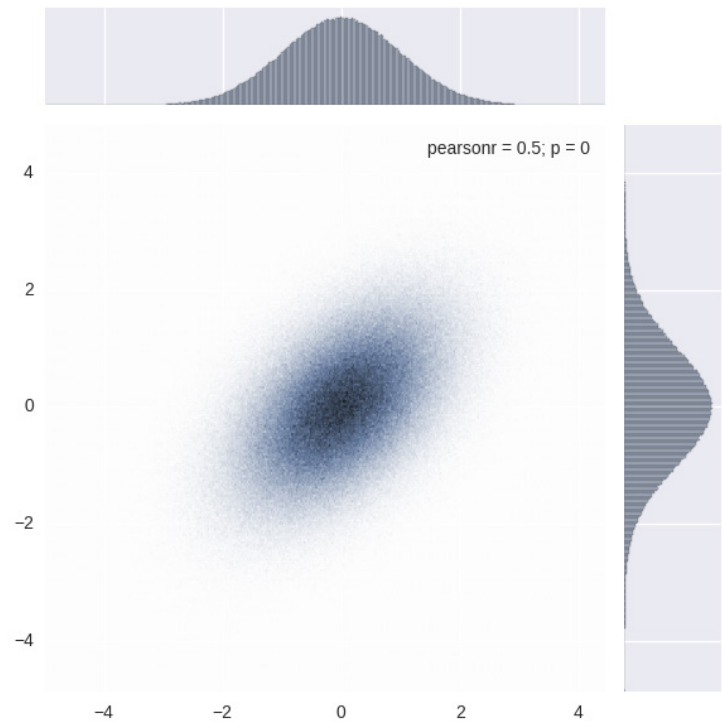
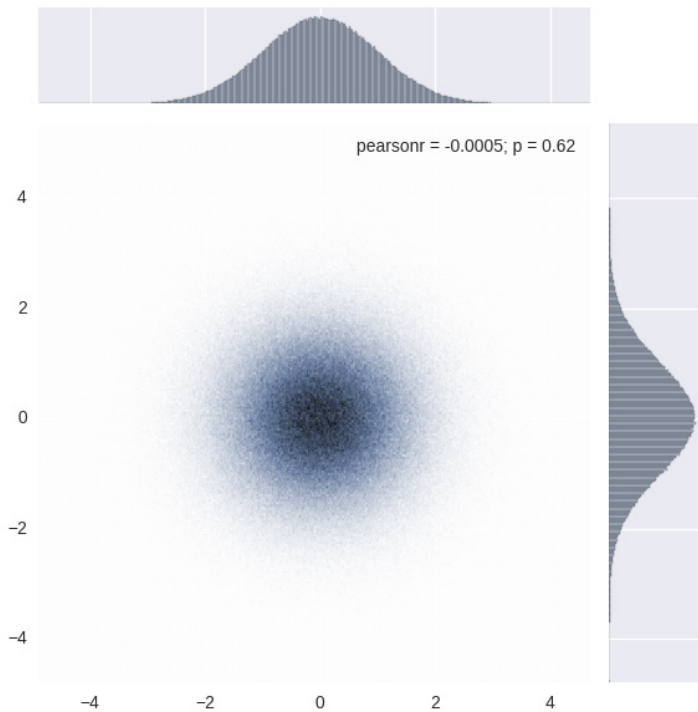
Covariance – 2D case



Adapted from Carl Henrik Ek's lecture (DD2434, 2015)

- Introduction
- Probabilistic approach
- **Probability basics**

Covariance – 2D case



Adapted from Carl Henrik Ek's lecture (DD2434, 2015)

- Introduction
- Probabilistic approach
- **Probability basics**

Central Limit Theorem

"The sum of a number of independent and identically distributed random variables with finite variances will tend to a normal distribution as the number of variables grows."

$$\mu = \mathbb{E}[X_i]$$
$$\sigma^2 = \text{var}[X_i]$$

$$M_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathcal{N}\left(\mu, \sigma / \sqrt{n}\right)$$

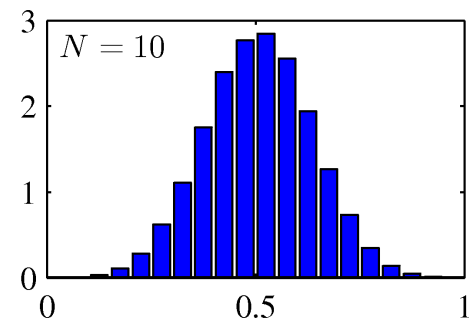
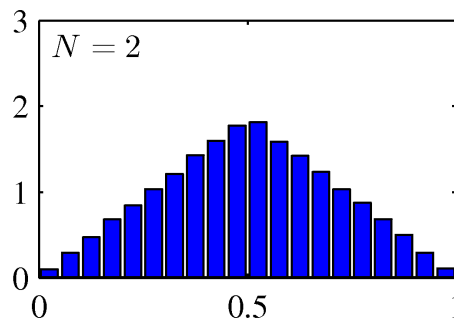
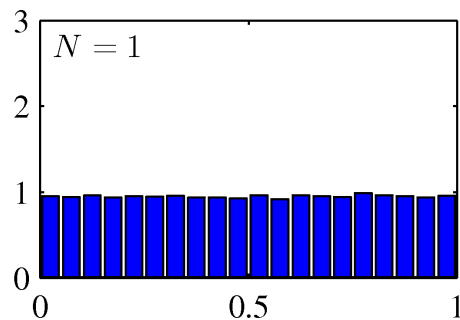
- Introduction
- Probabilistic approach
- **Probability basics**

Central Limit Theorem

"The sum of a number of independent and identically distributed random variables with finite variances will tend to a normal distribution as the number of variables grows."

$$\mu = \mathbb{E}[X_i]$$
$$\sigma^2 = \text{var}[X_i]$$

$$M_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathcal{N}\left(\mu, \sigma / \sqrt{n}\right)$$



Bishop, sec. 2.3